

Policies for Scientific Integrity and Reproducibility: Data and Code Sharing

Victoria Stodden
Department of Statistics
Columbia University

AAAS Annual Meeting
The Digitization of Science: Reproducibility and Interdisciplinary Knowledge Transfer
Feb 19, 2011

A Crisis in Computational Science

- Computational methods becoming central to the scientific enterprise:
 - enormous, and increasing, amounts of data collection,
 - intellectual contributions now encoded in software,
 - typical scientific results rely on both data and code.
- Data and code typically not made available, rendering published results unverifiable, not reproducible.

➡ A Credibility Crisis

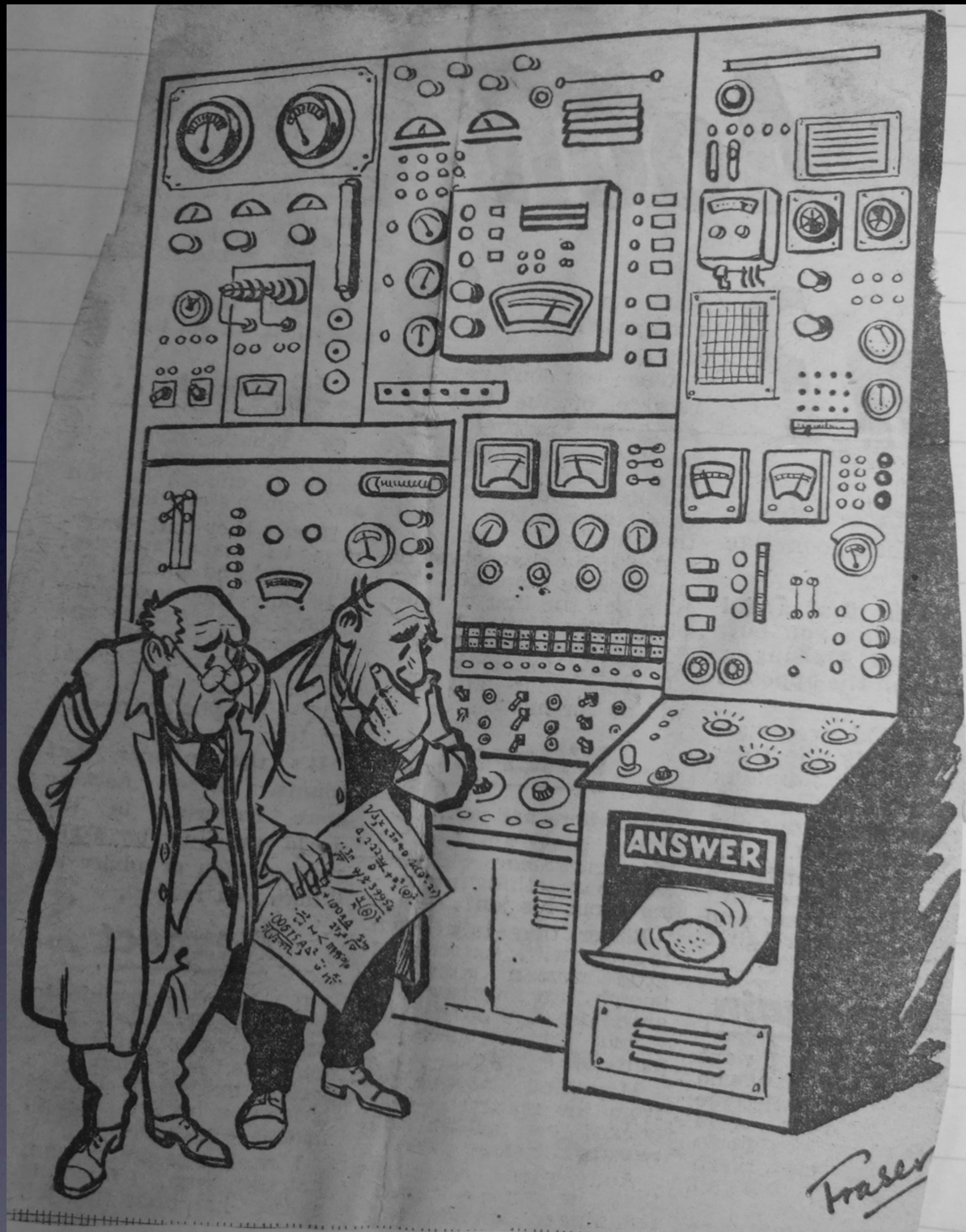
Computation Central to the Scientific Endeavor

For example, in statistics,

JASA June	Computational Articles	Code Publicly Available
1996	9 of 20	0%
2006	33 of 35	9%
2009	32 of 32	16%

Reproducibility is Central to Scientific Communication

- Other branches of science incorporate reproducibility of results:
 - deductive branch (mathematics, formal logic): the well-defined concept of the proof,
 - inductive branch (experimental sciences): machinery of hypothesis testing, structured communication of methods and protocols.
- Computational Science must develop standards for reproducibility before it can be considered a third branch of the scientific method,
➡ Data and Code Sharing, with publication.



Barriers to Reproducible Computational Science

Survey of Machine Learning Community (Stodden, 2010):

Code		Data
77%	Time to document and clean up	54%
52%	Dealing with questions from users	34%
44%	Not receiving attribution	42%
40%	Possibility of patents	-
34%	Legal Barriers (ie. copyright)	41%
-	Time to verify release with admin	38%
30%	Potential loss of future publications	35%
30%	Competitors may get an advantage	33%
20%	Web/disk space limitations	29%

Legal Barriers: Copyright

“To promote the Progress of Science and useful Arts, by securing for limited Times to Authors and Inventors the exclusive Right to their respective Writings and Discoveries.” (U.S. Const. art. I, §8, cl. 8)

- Original expression of ideas falls under copyright *by default* (papers, code, figures, tables..)
- Copyright secures exclusive rights vested in the author to:
 - reproduce the work
 - prepare derivative works based upon the original
 - limited time: generally life of the author +70 years

Exceptions and Limitations: Fair Use.

Responses Outside the Sciences I: Open Source Software

- Software with licenses that communicate alternative terms of use to code developers, rather than the copyright default.
- Hundreds of open source software licenses:
 - GNU Public License (GPL)
 - (Modified) BSD License
 - MIT License
 - Apache 2.0 License
 - ... see <http://www.opensource.org/licenses/alphabetical>

Responses Outside the Sciences 2: Creative Commons

- Adapts the Open Source Software approach to artistic and creative digital works
- Provides a suite of licensing options:
 - BY: if you use the work attribution must be provided,
 - NC: the work cannot be used for commercial purposes,
 - ND: no derivative works permitted,
 - SA: derivative works must carry the same license as the original

Response from Within the Sciences

The *Reproducible Research Standard (RRS)* (Stodden, 2009)

- A suite of license recommendations for computational science:
 - Release media components (text, figures) under CC BY,
 - Release code components under Modified BSD or similar,
 - Release data to public domain or attach attribution license.
- ➡ Remove copyright's barrier to reproducible research and,
- ➡ Realign the IP framework with longstanding scientific norms.

Winner of the Access to Knowledge Kaltura Award 2008

Benefits of the *RRS*

- Promotion of (legal) reproducible research,
- Focus becomes release of entire research compendium,
- Hook for funders, journals, institutional policy makers,
- Standardization avoids license incompatibilities,
- Clarity of rights, beyond Fair Use.

Yale Data and Code Sharing Roundtable 2009

- Roundtable on Data and Code Sharing in computational science
Nov 21, 2009:
 - gathered 30 computational scientists from a variety of fields, funding agency folks, publishers, librarians, university policy makers, lawyers...
 - Draft Position Statement (published in IEEE Computing in Science and Engineering, Sep/Oct 2010)
 - recommendations for stakeholders: scientists, journal editors, funding agencies, universities.
- <http://www.stanford.edu/~vcs/Conferences/RoundtableNov212009/>

References

- “Enabling Reproducible Research: Open Licensing for Scientific Innovation”
- “The Scientific Method in Practice: Reproducibility in the Computational Sciences”

<http://www.stanford.edu/~vcs>

- Data and Code Sharing Roundtable, Nov 2009
- Reproducible Research: Tools and Strategies for Scientific Computing, July 2011
- Reproducible Research in Computational Science: What, Why and How, Community Forum, July 2011